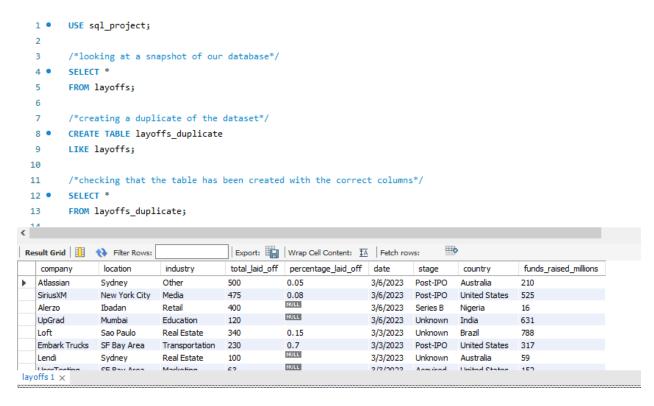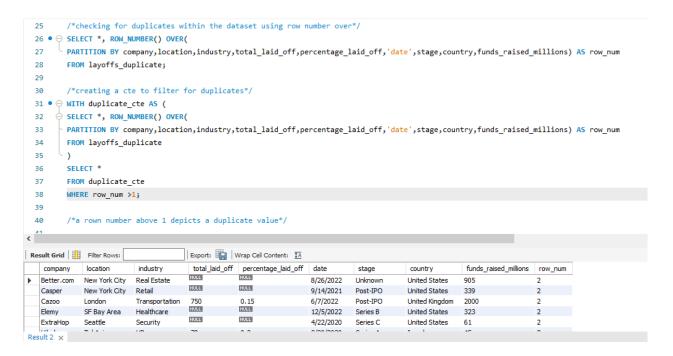In this project, I will clean up a database and perform exploratory data analysis on it. I will do this by

- removing duplicates
- standardizing data
- addressing null values or blanks
- removing any rows and columns not needed
- conduct an EDA

I start off by looking at a snapshot of our data. After first impressions, I create a duplicate of the database and store the original away to be referred to whenever needed.
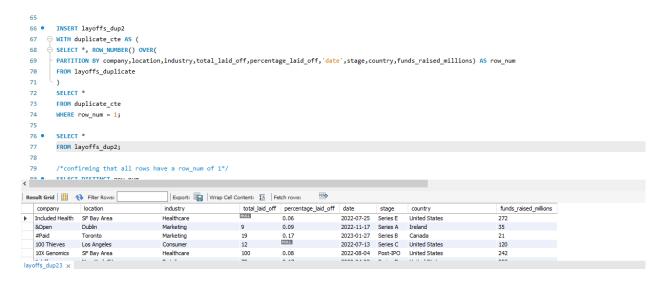
```
1   USE sql_project;
2
3   /*looking at a snapshot of our database*/
4   SELECT *
5   FROM layoffs;
6
7   /*creating a duplicate of the dataset*/
8   CREATE TABLE layoffs_duplicate
9   LIKE layoffs;
10
11  /*checking that the table has been created with the correct columns*/
12  SELECT *
13  FROM layoffs_duplicate;
14
```

| company | location | industry | total_laid_off | percentage_laid_off | date | stage | country | funds_raised_millions |
|---|---|---|---|---|---|---|---|---|
| Atlassian | Sydney | Other | 500 | 0.05 | 3/6/2023 | Post-IPO | Australia | 210 |
| SiriusXM | New York City | Media | 475 | 0.08 | 3/6/2023 | Post-IPO | United States | 525 |
| Alerzo | Ibadan | Retail | 400 | NULL | 3/6/2023 | Series B | Nigeria | 16 |
| UpGrad | Mumbai | Education | 120 | NULL | 3/6/2023 | Unknown | India | 631 |
| Loft | Sao Paulo | Real Estate | 340 | 0.15 | 3/3/2023 | Unknown | Brazil | 788 |
| Embark Trucks | SF Bay Area | Transportation | 230 | 0.7 | 3/3/2023 | Post-IPO | United States | 317 |
| Lendi | Sydney | Real Estate | 100 | NULL | 3/3/2023 | Unknown | Australia | 59 |
| UserTesting | SF Bay Area | Marketing | 63 | NULL | 3/3/2023 | Acquired | United States | 152 |

layoffs 1 ✕

Using row_number() over(), I am able to suss out where a particular row has repeat datapoints in all columns meaning it is a duplicate.

```
25     /*checking for duplicates within the dataset using row number over*/
26  ●  ⊖ SELECT *, ROW_NUMBER() OVER(
27        PARTITION BY company,location,industry,total_laid_off,percentage_laid_off,'date',stage,country,funds_raised_millions) AS row_num
28        FROM layoffs_duplicate;
29
30        /*creating a cte to filter for duplicates*/
31  ●  ⊖ WITH duplicate_cte AS (
32     ⊖ SELECT *, ROW_NUMBER() OVER(
33        PARTITION BY company,location,industry,total_laid_off,percentage_laid_off,'date',stage,country,funds_raised_millions) AS row_num
34        FROM layoffs_duplicate
35        )
36        SELECT *
37        FROM duplicate_cte
38        WHERE row_num >1;
39
40        /*a rown number above 1 depicts a duplicate value*/
```

| company | location | industry | total_laid_off | percentage_laid_off | date | stage | country | funds_raised_millions | row_num |
|---|---|---|---|---|---|---|---|---|---|
| Better.com | New York City | Real Estate | NULL | NULL | 8/26/2022 | Unknown | United States | 905 | 2 |
| Casper | New York City | Retail | NULL | NULL | 9/14/2021 | Post-IPO | United States | 339 | 2 |
| Cazoo | London | Transportation | 750 | 0.15 | 6/7/2022 | Post-IPO | United Kingdom | 2000 | 2 |
| Elemy | SF Bay Area | Healthcare | NULL | NULL | 12/5/2022 | Series B | United States | 323 | 2 |
| ExtraHop | Seattle | Security | NULL | NULL | 4/22/2020 | Series C | United States | 61 | 2 |

Result 2 ✕

I create a second duplicate table and populate it with the unique data points, removing the duplicates.

```
52      /*creating a new table that contains only unique rows*/
53  ●  ⊖ CREATE TABLE layoffs_dup2(
54        company TEXT,
55        location TEXT,
56        industry TEXT,
57        total_laid_off INT,
58        percentage_laid_off TEXT,
59        date TEXT,
60        stage TEXT,
61        country TEXT,
62        funds_raised_millions INT,
63        row_num INT
64        );
65
```

```
65
66 ●    INSERT layoffs_dup2
67  ⊝   WITH duplicate_cte AS (
68  ⊝   SELECT *, ROW_NUMBER() OVER(
69    │  PARTITION BY company,location,industry,total_laid_off,percentage_laid_off,'date',stage,country,funds_raised_millions) AS row_num
70    │  FROM layoffs_duplicate
71    ⊢  )
72       SELECT *
73       FROM duplicate_cte
74       WHERE row_num = 1;
75
76 ●    SELECT *
77       FROM layoffs_dup2;
78
79       /*confirming that all rows have a row_num of 1*/
80 ●    SELECT DISTINCT row num
```

| company | location | industry | total_laid_off | percentage_laid_off | date | stage | country | funds_raised_millions |
|---|---|---|---|---|---|---|---|---|
| Included Health | SF Bay Area | Healthcare | NULL | 0.06 | 2022-07-25 | Series E | United States | 272 |
| &Open | Dublin | Marketing | 9 | 0.09 | 2022-11-17 | Series A | Ireland | 35 |
| #Paid | Toronto | Marketing | 19 | 0.17 | 2023-01-27 | Series B | Canada | 21 |
| 100 Thieves | Los Angeles | Consumer | 12 | NULL | 2022-07-13 | Series C | United States | 120 |
| 10X Genomics | SF Bay Area | Healthcare | 100 | 0.08 | 2022-08-04 | Post-IPO | United States | 242 |

layoffs_dup23 ✕

We proceed using the second duplicate table as our working dataset.

On to standardization, the column "company" contains words and spaces. The spaces sometimes come before the words or after, creating inconsistencies that may affect future analysis. I us the trim function to solve this issue.

```
73       FROM duplicate_cte
74       WHERE row_num = 1;
75
76 ●    SELECT *
77       FROM layoffs_dup2;
78
79       /*confirming that all rows have a row_num of 1*/
80 ●    SELECT DISTINCT row_num
81       FROM layoffs_dup2;
82
83       /*standardizing the data*/
84
85       /*removing any white spaces surrounding company*/
86 ●    UPDATE layoffs_dup2
87       SET company = TRIM(company);
88
89 ●    SELECT DISTINCT industry
90       FROM layoffs_dup2
91       ORDER BY 1;
92
93  ⊝   /*we notice a few issues in the industry column to be addressed
94    │  -cryto / crypto currency
95    │  -null values
96    ⊢  -blanks*/
```

I move to the next column "industry". The industry column has entries under both "Crypto" and "Cryto Currency" which I decide are the same thing. Since "Cryto" appears more frequently, I would be using that. Industry also has some null values and blanks.

```
88
89 ●    SELECT DISTINCT industry
90       FROM layoffs_dup2
91       ORDER BY 1;
92
93   ⊖  /*we notice a few issues in the industry column to be addressed
94    │   -cryto / crypto currency
95    │   -null values
96    └   -blanks*/
97
98 ●    SELECT *
99       FROM layoffs_dup2
100      WHERE industry LIKE 'Crypto%';
101
102 ●   UPDATE layoffs_dup2
103      SET industry = 'Crypto'
104      WHERE industry LIKE 'Crypto%';
105
106 ●   SELECT DISTINCT location
107      FROM layoffs_dup2
108      ORDER BY 1;
109
110 ●   SELECT DISTINCT country
111      FROM layoffs_dup2
```

Next, I look at where a company in the same location appears more than once but has industry nulls or blanks in one of the entries. To make things easier, I will set the blanks to null values, the populate the nulls with industry values inferred from the existing data.

```
/*addressing null values*/
/*where do we see the same company in the same location having industry null and not null*/
SELECT *
FROM layoffs_dup2 AS l1
JOIN layoffs_dup2 AS l2
USING (company)
WHERE (l1.industry IS NULL OR l1.industry = '')
AND l2.industry IS NOT NULL;

/*setting blank spaces to null*/
UPDATE layoffs_dup2
SET industry = NULL
WHERE industry = '';

/*inputing inferred value for industry column from existing data*/
UPDATE layoffs_dup2 AS l1
JOIN layoffs_dup2 AS l2
    USING (company)
SET l1.industry = l2.industry
WHERE l1.industry IS NULL
AND l2.industry IS NOT NULL;

/*checking that inferable industry null values have been populated*/
SELECT *
```

After confirming that inferable entries have been made, I move to the next column "country". There are 2 different "United States" because one has a period at the end. We need to remove said period so that the system recognizes them as one.

```
128      /*checking that inferable industry null values have been populated*/
129 •    SELECT *
130      FROM layoffs_dup2
131      WHERE industry IS NULL;
132
133 •    SELECT DISTINCT location
134      FROM layoffs_dup2
135      ORDER BY 1;
136
137 •    SELECT DISTINCT country
138      FROM layoffs_dup2
139      ORDER BY 1;
140      /*united states is showing a little discrepancy*/
141 •    SELECT *
142      FROM layoffs_dup2
143      WHERE country LIKE 'United States%';
144
145 •    UPDATE layoffs_dup2
146      SET country = TRIM(TRAILING '.' FROM country)
147      WHERE country LIKE 'United States%';
148
149 •    SELECT DISTINCT country
150      FROM layoffs_dup2
151      ORDER BY 1;
```

Checking the data types of the columns, I discover that date column is a text data type. I convert it to the more appropriate 'date' data type.

```
152
153      /*checking the datatypes of our columns*/
154 •    DESCRIBE layoffs_dup2;
155
156      /*converting the date to date format*/
157 •    UPDATE layoffs_dup2
158      SET date = STR_TO_DATE(date, '%m/%d/%Y');
159
160      /*now that it's in a date format, we can easily convert to date datatype*/
161 •    ALTER TABLE layoffs_dup2
162      MODIFY COLUMN date DATE;
163
```

My analysis will focus on the numeric aspect of the dataset, having null values the numeric columns will be problematic.

```
SELECT *
FROM layoffs_dup2
WHERE total_laid_off IS NULL
AND percentage_laid_off IS NULL;
/*348 rows have both these key columns null, roughly 15% of our dataset. seeing no viable way to populate
these cells and yet realising not removing them will heavily influence my analysis moving forward. I opt to delete
them. I can always fall back on the original database if need be*/

DELETE
FROM layoffs_dup2
WHERE total_laid_off IS NULL
AND percentage_laid_off IS NULL;

/*row_num column is now redundant*/
ALTER TABLE layoffs_dup2
DROP COLUMN row_num;
```

Checking for null values in our 2 numeric columns, I find that 348 rows have both numeric columns null, this accounts for about 15% of our original data. Unfortunately, I can not infer these values from the existing data, nor can I confidently input a calculated estimate. I decided that leaving these rows in will affect my analysis so opt to remove them. Also, the row_num column I used to find the duplicates is no longer needed so I remove that as well.

Feeling content with the cleaning for now, I move on to my exploratory data analysis. Things I am curious about the data set include:

- What is the timeframe of the data set ?
- Which companies went completely under during the period ?
- Which companies are laying off the most employees ?
- Which industry laid off the most on average ?
- Which countries laid off the most ?
- How did layoffs progress over the period of interest ?
- How the stage a company is in affected layoffs.

```sql
/*what date range are we working with in this database*/
SELECT MIN(date),MAX(date)
FROM layoffs_dup2;

/*which companies went completely under */
SELECT *
FROM layoffs_dup2
WHERE percentage_laid_off = 1
ORDER BY funds_raised_millions DESC;

/*which comapnies are laying off a large number of employees*/
SELECT company,SUM(total_laid_off)
FROM layoffs_dup2
GROUP BY company
ORDER BY 2 DESC;

/*what industry got hit the most*/
SELECT industry,AVG(total_laid_off)
FROM layoffs_dup2
GROUP BY industry
ORDER BY 2 DESC;
```

The data set is from 11th March 2020 to to 6th March 2023, starting at the peak of the world wide outbreak of the corona virus and afterwards.

```
187
188     /*which companies went completely under */
189 •   SELECT industry, COUNT(company)
190     FROM layoffs_dup2
191     WHERE percentage_laid_off = 1
192     GROUP BY industry
193     ORDER BY 2 DESC;
```

| industry | COUNT(company) |
|---|---|
| Retail | 13 |
| Food | 13 |
| Finance | 12 |
| Education | 9 |
| Healthcare | 7 |
| Crypto | 7 |
| Transportation | 6 |
| Real Estate | 6 |
| Marketing | 5 |
| Media | 5 |
| Consumer | 5 |
| Travel | 4 |
| Other | 4 |
| Infrastructure | 3 |
| Product | 3 |
| Recruiting | 2 |
| Support | 2 |
| Fitness | 2 |

Result 5 ✕

At the top of the list of industries we have retail, food and finance.

```
194
195       /*which comapnies are laying off a large number of employees*/
196  ●   SELECT company,SUM(total_laid_off)
197       FROM layoffs_dup2
198       GROUP BY company
199       ORDER BY 2 DESC;
```

| company | SUM(total_laid_off) |
| --- | --- |
| Amazon | 18150 |
| Google | 12000 |
| Meta | 11000 |
| Salesforce | 10090 |
| Microsoft | 10000 |
| Philips | 10000 |
| Ericsson | 8500 |
| Uber | 7585 |
| Dell | 6650 |
| Booking.com | 4601 |
| Cisco | 4100 |
| Peloton | 4084 |
| Byju's | 4000 |
| Carvana | 4000 |
| Twitter | 3940 |
| Better.com | 3900 |
| IBM | 3900 |
| Groupon | 3800 |
| Bytedance | 3750 |

Result 6 ×

Large companies like Amazon, Google and Meta laid off the most employees during the period. However, these companies also have a disproportionately large number of employees in general.

```
200
201     /*what industry got hit the most*/
202 ●   SELECT industry,COUNT(company) ,SUM(total_laid_off)
203     FROM layoffs_dup2
204     GROUP BY industry
205     ORDER BY 3 DESC;
206
207     /*what countries got hit the most*/
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| industry | COUNT(company) | SUM(total_laid_off) |
|---|---|---|
| Consumer | 101 | 45182 |
| Retail | 163 | 43613 |
| Other | 106 | 36209 |
| Transportation | 128 | 33548 |
| Finance | 239 | 28344 |
| Healthcare | 163 | 25894 |
| Food | 115 | 22855 |
| Real Estate | 100 | 17565 |
| Travel | 56 | 17159 |
| Hardware | 13 | 13828 |
| Education | 79 | 13338 |
| Sales | 34 | 13216 |
| Crypto | 89 | 10693 |
| Marketing | 123 | 10258 |
| Fitness | 26 | 8748 |
| Security | 62 | 5979 |
| Infrastructure | 32 | 5785 |

Consumer and retail are leading industries in layoffs.

```
206
207     /*what countries got hit the most*/
208 ●   SELECT country,COUNT(country),SUM(total_laid_off)
209     FROM layoffs_dup2
210     GROUP BY country
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: IA

| country | COUNT(country) | SUM(total_laid_off) |
|---|---|---|
| United States | 1291 | 256420 |
| India | 139 | 35793 |
| Canada | 90 | 6319 |
| Brazil | 69 | 10391 |
| Germany | 62 | 8701 |
| United Kingdom | 58 | 6398 |
| Israel | 47 | 3638 |
| Australia | 46 | 2324 |
| Singapore | 26 | 5995 |
| Indonesia | 20 | 3521 |
| Sweden | 17 | 11264 |
| China | 13 | 5905 |
| Netherlands | 12 | 17220 |
| Nigeria | 11 | 1882 |
| United Arab Emirates | 6 | 995 |
| Kenya | 6 | 349 |
| France | 5 | 915 |
| Argentina | 5 | 323 |
| Estonia | 5 | 333 |
| Ireland | 4 | 257 |

USA contributing the most companies to the database naturally also has the most number of
layoffs.

```
214 •  SELECT YEAR(date),COUNT(DISTINCT(MONTH(date))),SUM(total_laid_off)
215     FROM layoffs_dup2
216     WHERE date IS NOT NULL
217     GROUP BY YEAR(date)
218     ORDER BY 1;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: ‡A

| YEAR(date) | COUNT(DISTINCT(MONTH(date))) | SUM(total_laid_off) |
|---|---|---|
| 2020 | 10 | 80998 |
| 2021 | 11 | 15823 |
| 2022 | 12 | 160322 |
| 2023 | 3 | 125677 |

The layoffs were lowest in 2020, but subsequent years saw higher numbers in the aftermath of
the pandemic. Though 2023 shows a reduction in layoffs, only 3 months out of the year are
accounted for.

```
232     /*over the period, which companies have layed the most employees off*/
233 •  WITH company_year(company, year,total_layoffs) AS
234   ⊖ (
235     SELECT company, YEAR(date),SUM(total_laid_off)
236     FROM layoffs_dup2
237     GROUP BY company,YEAR(date)
238     ORDER BY company ASC),
239   ⊖ company_year_rank AS (SELECT *,
240     DENSE_RANK() OVER(PARTITION BY year ORDER BY total_layoffs DESC) AS ranking
241     FROM company_year
242     WHERE year IS NOT NULL)
243     SELECT *
244     FROM company_year_rank
245     WHERE ranking <= 5
246     ORDER BY year;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: ‡A

| company | year | total_layoffs | ranking |
|---|---|---|---|
| Uber | 2020 | 7525 | 1 |
| Booking.com | 2020 | 4375 | 2 |
| Groupon | 2020 | 2800 | 3 |
| Swiggy | 2020 | 2250 | 4 |
| Airbnb | 2020 | 1900 | 5 |

Result 17 ✕

The top 5 companies laying off employees over the period was quite volatile.

```
248     /*does the stage a company is in affect layoffs*/
249 ●   SELECT stage,SUM(total_laid_off)
250     FROM layoffs_dup2
251     GROUP BY stage
252     ORDER BY 2 DESC;
```

**Result Grid** | Filter Rows: | Export: | Wrap Cell Content: ⊺A

| stage | SUM(total_laid_off) |
|---|---|
| Post-IPO | 204073 |
| Unknown | 40716 |
| Acquired | 27496 |
| Series C | 20017 |
| Series D | 19225 |
| Series B | 15311 |
| Series E | 12697 |
| Series F | 9932 |
| Private Equity | 7957 |
| Series H | 7244 |
| Series A | 5678 |
| Series G | 3697 |
| Series J | 3370 |
| Series I | 2855 |
| Seed | 1636 |
| Subsidiary | 1094 |
| NULL | 322 |

Result 18 ✕

Finally, companies in Post-IPO stage laid off the most during the period.